

# Complete genomic characterization of two *Escherichia coli* lineages responsible for a cluster of carbapenem-resistant infections in a Chinese hospital

Zong, Zhiyong; Fenn, Samuel; Connor, Christopher; Feng, Yu; McNally, Alan

DOI:

[10.1093/jac/dky210](https://doi.org/10.1093/jac/dky210)

License:

Other (please specify with Rights Statement)

*Document Version*

Peer reviewed version

*Citation for published version (Harvard):*

Zong, Z, Fenn, S, Connor, C, Feng, Y & McNally, A 2018, 'Complete genomic characterization of two *Escherichia coli* lineages responsible for a cluster of carbapenem-resistant infections in a Chinese hospital', *Journal of Antimicrobial Chemotherapy*. <https://doi.org/10.1093/jac/dky210>

[Link to publication on Research at Birmingham portal](#)

## **Publisher Rights Statement:**

This is a pre-copyedited, author-produced PDF of an article accepted for publication in *Journal of Antimicrobial Chemotherapy* following peer review. The version of record Zhiyong Zong, Samuel Fenn, Christopher Connor, Yu Feng, Alan McNally; Complete genomic characterization of two *Escherichia coli* lineages responsible for a cluster of carbapenem-resistant infections in a Chinese hospital, *Journal of Antimicrobial Chemotherapy*, , dky210, <https://doi.org/10.1093/jac/dky210> is available online at: 10.1093/jac/dky210

## **General rights**

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

## **Take down policy**

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.

**Complete genomic characterisation of two *Escherichia coli* lineages responsible for a cluster of carbapenem resistant infections in a Chinese hospital**

Zhiyong ZONG, Samuel FENN, Christopher CONNOR, Yu FENG, Alan McNALLY\*

<sup>1</sup>Centre for Infectious Diseases, West China Hospital of Sichuan University, Chengdu, China

<sup>2</sup>Institute of Microbiology and Infection, College of Medical and Dental Science, University of Birmingham, Birmingham, United Kingdom, B15 2TT

\*Corresponding author: Dr Alan McNally, Institute of Microbiology and Infection, College of Medical and Dental Science, University of Birmingham, Birmingham B15 2TT. 0044 121 4158433. a.mcnally.1@bham.ac.uk

Running title: Carbapenem resistant clones of *E. coli*

## Abstract

Objectives: The increase in infections as a result of multi-drug resistant strains of *Escherichia coli* is a global health crisis. The emergence of globally disseminated lineages of *E. coli* carrying ESBL genes has been well characterised. An increase in strains producing carbapenemase enzymes and mobile colistin resistance is now being reported, but to date there is little genomic characterisation of such strains.

Methods: Routine screening of patients within an ICU of West China Hospital identified a number of *E. coli* carrying the *bla*<sub>NDM-5</sub> carbapenemase gene, found to be two distinct clones, *E. coli* ST167 and ST617.

Results: Interrogation of publically available data shows isolation of ESBL and carbapenem resistant strains of both lineages from clinical cases across the world. Further analysis of a large collection of publically available genomes shows that ST167 and ST617 have emerged in distinct patterns from the ST10 clonal complex of *E. coli*, but share evolutionary events involving switches in LPS genetics, intergenic regions and anaerobic metabolism loci.

Conclusions: The identification of these lineages of *E. coli* and their shared genetic traits suggest there may be evolutionary events which underpin the emergence of carbapenem resistance plasmid carriage in *E. coli*.

## Introduction

Infections from multi-drug resistant (MDR) *Escherichia coli* are a significant global health care threat.<sup>1</sup> MDR in *E. coli* is largely confined to strains capable of causing extra-intestinal infections (ExPEC) such as urinary tract infections (UTI) and bacteraemia.<sup>1-4</sup> As many as 50% of *E. coli* strains isolated from UTI and bacteraemia cases may exhibit resistance to three or more classes of antibiotic, termed MDR. This resistance is primarily driven by the acquisition of large plasmids containing multiple resistance genes.<sup>2</sup> The rapid global dissemination of MDR *E. coli* is associated with carriage of plasmids containing genes encoding extended-spectrum  $\beta$ -lactamases (ESBL) which confer resistance to third-generation cephalosporins.<sup>5</sup> The carriage of MDR plasmids containing ESBL genes renders *E. coli* susceptible only to the carbapenem class of antibiotics and the antimicrobial compound colistin.<sup>5</sup> However strains of *E. coli* are now being reported with plasmids containing  $\beta$ -lactamases conferring resistance to carbapenems (carbapenemases) and the *mcr-1* colistin resistance gene.<sup>6-9</sup>

The global dissemination of ESBL *E. coli* is attributable to the rapid dispersal of a small number of *E. coli* lineages. The most dominant of these is the ST131 lineage which is predominantly associated with carriage of the *bla*<sub>CTX-M-15</sub> ESBL gene.<sup>2</sup> ST131 is an ExPEC lineage and the most common cause of UTI and bacteraemia in the developed world.<sup>2</sup> Other dominant lineages of ESBL *E. coli* are ST73, ST95, and ST648 which are also ExPEC.<sup>3,4</sup> ESBL carriage can also be found transiently in strains belonging the ST10 clonal complex of *E. coli*.<sup>3</sup> ST10 complex strains are host generalist *E. coli* which are frequently found as intestinal commensal inhabitants of mammals and avian species,<sup>10</sup> and are devoid of the virulence-associated genes known to be required for pathogenesis.<sup>11</sup> Our knowledge of the genomic landscape

of carbapenemase production in *E. coli* is far less developed, with the vast majority of reports being genomes of individual clinical isolates sporadically distributed across the globe. Just one significant publication exists reporting a specifically designed genomic analysis of a temporal collection of carbapenem resistant *E. coli* which showed very wide dissemination of carbapenem resistance across species and within-species lineages of the enterobacteriaceae.<sup>12</sup>

Here we report the isolation of *E. coli* containing the carbapenem-resistance gene *bla*<sub>NDM-5</sub> in an ICU ward in West China Hospital, Chengdu. The isolates do not belong to one of the dominant MDR lineages of ExPEC, but to ST167 and ST617, both members of the ST10 clonal complex. Genomic data supports the long-term presence of these bacteria in the ICU with repeated dissemination from a central reservoir. Contextualisation of the Chinese strains with a collection of publically available genomes shows isolation of MDR ST167 and ST617 strains from clinical episodes across the world, and in the case of ST167 frequent occurrence of carriage of both ESBL and carbapenemase genes. By comparing these lineages to a large number of publically available ST10 genomes we identify potentially significant events in their evolutionary trajectories, including mutations in the LPS biosynthesis locus which truncate LPS. We also find evidence of compensatory mutations in intergenic regions as found in *E. coli* ST131 as well as mutations in anaerobic metabolism loci. Our findings support the need for a more concerted global surveillance effort focussing on identifying frequently occurring lineages of carbapenem resistant *E. coli*.

## Methods

### Bacterial isolation and characterisation

Strain 0215 was recovered from a rectal swab of a 75-year-old male patient on September 2013 in a 50-bed medical ICU at West China Hospital, Chengdu, during routine screening that is performed as standard in the ICU on all new admissions. During a 7-month period from May to November 2014, *bla*<sub>NDM-5</sub> positive *E. coli* were recovered from the rectal swabs of 8 different patients (Supplementary Table S1) from a total of 560 patients admitted to the ICU during this period. Furthermore, one of the 8 patients developed bacteraemia during his ICU stay and an *E. coli* was recovered from his blood and included in the study. During the study period, two additional *E. coli* clinical isolates carrying *bla*<sub>NDM-5</sub> were recovered in the hospital, from two patients on admission. Rectal swabs were collected from patients within 2 days of admission to the ICU and within the 3 days prior to ICU discharge for those patients with a length of stay of 3 days or more. Swabs were transferred to the laboratory in transport media and were screened for carbapenem-resistant Enterobacteriaceae using the CHROMAgar Orientation agar plates containing 2 µg/mL meropenem.

## **Ethics**

This study was conducted in accordance with the amended Declaration of Helsinki and was approved, under a waiver of consent, by the Ethics Committee of West China Hospital. Rectal swabs were collected from patients within 2 days of admission to the ICU and within the 3 days prior to ICU discharge for those patients with a length of stay of 3 days or more.

## **Genome sequencing**

The ST167 and ST617 strains isolated in Chengdu were cultured in LB broth at 37°C overnight. DNA was extracted using QIAamp<sup>®</sup> DNA Mini Kit (QIAGEN) and 150 bp paired-end libraries of each strain prepared and sequenced using the Illumina HiSeq

X-Ten platform (raw data accession numbers Table S2 and S3). Genomes were assembled using SPAdes<sup>13</sup> and annotated using Prokka.<sup>14</sup> The MLST sequence type of the strains was determined using the in silico prediction tool MLSTFinder.<sup>15</sup> The *E. coli* genome database Enterobase ([www.enterobase.warwick.ac.uk](http://www.enterobase.warwick.ac.uk)) was interrogated on 1<sup>st</sup> December 2016 and all available ST167 and ST617 genomes were downloaded (Table S2 and S3) and annotated using Prokka. A further 256 ST10 genomes were selected to represent the geographical, temporal, and source attribution diversity present in the database (Table S4) and were downloaded and annotated using Prokka. To select these genomes a phylogenetic tree was inferred from the assembled genome of every ST10 on Enterobase using Parsnp.<sup>16</sup> From this phylogeny 500 genomes were chosen to span the entire phylogenetic diversity, and then the final selection made to represent the full ST10 diversity as described. The antibiotic resistance gene profile of all isolates was determined using Abricate (<https://github.com/tseemann/abrigate>).

#### **High-resolution SNP analysis**

We created a closed genome sequence for a Chinese ST167 strain 1237 by combining our Illumina sequence data with data generated on the Minlon sequencer. Raw Minlon reads were converted into fastQ format (accession number PRJNA422975) using Poretools<sup>17</sup> and assembled using Canu,<sup>18</sup> resulting in a single contig chromosome and four distinct single contig plasmids. The raw illumina data was then used to polish the genome assembly via five iterative rounds of polishing with Pilon.<sup>19</sup> The ST167 and ST617 genomes from Chengdu were analysed by mapping raw reads against the hybrid assembled ST167 genome. Mapping was performed using Snippy (<https://github.com/tseemann/snippy>) and the resulting SNP profiles were used to create a consensus sequence for each genome

which was aligned using the parsnp alignment tool in Harvest.<sup>16</sup> Analysis of the plasmid containing the *bla*<sub>NDM-5</sub> gene revealed that it was a 47-kb IncX3 plasmid and there were no antibiotic resistant genes other than *bla*<sub>NDM-5</sub> located on the plasmid. Specific mapping of the raw Illumina data against the pNDM5 plasmid was performed for all strains as described above.

#### **Phylogenetic analysis**

Pan-genomes were constructed for the ST167, ST617, ST10, and combined datasets using Roary<sup>20</sup> with the --e --mafft setting to create a concatenated alignment of core CDS. The alignments were used to infer ST167, ST617, ST10, and combined phylogenies using RaxML<sup>21</sup> with the GTR-Gamma model of site heterogeneity and 100 bootstrap iterations. Carriage of ESBL and carbapenemase genes was annotated on the trees using Phandango (<https://jameshadfield.github.io/phandango/>), and geographical source was annotated using iTOL.<sup>22</sup>

#### **Detection of lineage specific genetic traits**

Microbial GWAS was performed using two approaches. First the combined data set pan-genome matrix was used as input for Scoary<sup>23</sup> searching for loci unique to ST167, ST617, and both ST167 and ST617 versus ST10. In parallel we also used SEER<sup>24</sup> to detect kmers significantly associated with ST167, ST617, or both combined versus ST10. The results of both approaches were combined to identify coding loci associated with the emergence of ST167 and ST617. In silico serotyping was performed using two independent methods, SRST2 and SerotypeFinder.<sup>25,26</sup> Both methods utilise WGS data to specific O and H antigens to strains. Intergenic regions (IGRs) were investigated using Piggy<sup>27</sup> to search for IGRs which had switched<sup>28</sup> in ST617, ST167, or both compared to ST10. This data was combined



with SEER data to identify high-confidence IGR switches associated with the emergence of ST167 and ST617.

## Results

### **Presence of *E. coli* ST167 and ST617 strains containing the NDM-5 carbapenemase resistance gene in an ICU ward in West China Hospital.**

A total of ten isolates of *E. coli* containing *bla*<sub>NDM-5</sub> were obtained during the investigation. Nine of these isolates belonged to sequence types ST167/617 (Table S1), which are members of the ST10 complex of *E. coli* most commonly associated with mammalian intestinal commensal carriage. Three ST167 isolates (0215, 243 and 25) were obtained from swabs or clinical samples collected on admission to hospital, suggesting that they were introduced from external sources. The three patients were all citizens of Chengdu city but they were admitted to different local hospitals before transferring to West China hospital. The remaining ST167 isolates were recovered from swabs or samples collected at least 3 days after admission to the ICU of West China hospital, from patients whose initial swabs were CRE negative, indicating that they were acquired during their ICU stay. ST167 *E. coli* carrying *bla*<sub>NDM-5</sub> caused infections (bacteremia and abdominal infection) in only two patients but colonised the others. Both ST617 *E. coli* carrying *bla*<sub>NDM-5</sub> only colonised patients. All patients colonised or infected with *E. coli* carrying *bla*<sub>NDM-5</sub> of ST167 or ST617 had received carbapenems before the recovery of the isolates.

### **SNP analysis suggests continued dissemination of strains from a central reservoir and sharing of resistance plasmid between lineages.**

To determine the level of relatedness between all isolated strains we mapped reads of all the strains against a closed ST167 strain (strain 1237) generated by a combination of Illumina and Minlon sequence data. The resulting high-resolution

SNP alignment showed the distance between the ST167 and ST617 strains to be over 25,000 SNPs, confirming they are distinct lineages, with the two ST617 isolates separated by just 7 SNPs. Deeper analysis of the ST167 cluster of strains showed diversity ranging from 5 to 799 SNPs (Fig 1). Strains 936 and 1222 (both carriage isolates) are the most closely related isolates with just 5 SNPs difference between them, with both strains being acquired by patients in the ICU within one month of each other. However these strains are 73 SNPs different from a strain isolated the exact same month on the ICU from a strain (1237) that was acquired in the ICU. This is almost double the genetic distance (46 SNPs) from a strain acquired (442 and 57, isolated from the same patient) in the ICU two months earlier. These distances are also larger than those for any isolate to the first two strains brought into the ICU, strain 0215 and strain 243, which differ from all other isolates by around 30 SNPs, and from each other by 15 SNPs. Such an observation suggests a potential combination of patient-to-patient transmission in the affected ICU,<sup>29</sup> along with the continued dissemination of the strain from a central reservoir where there is an accumulation of diversity.<sup>29,30</sup> Genomic analysis also allows us to identify a second introgression of an ST167 strain (25) from the community, which is over 700 SNPs different from the other isolates. Mapping of the raw sequence data against the 43kb IncX3 plasmid containing *bla*<sub>NDM-5</sub> also confirmed that the plasmid present in the ST617 strains was identical to that in all of the ST167 strains with just two detectable SNPs difference across the isolates.

#### **MDR ST167 and ST617 *E. coli* have been isolated across the world.**

We sought to contextualise the wider relevance of our Chengdu isolates by investigating the wider prevalence of ST167 and ST617 strains. We searched the Enterobase *E. coli* database and recovered a total of 87 genomes of ST167 (table

S2) and 86 genomes of ST617 (table S3), isolated from across the world. A core CDS-based phylogeny of both lineages showed a diverse set of genomes with around 17,000 SNPs in ST167 and around 15,000 SNPs in ST617. Annotation of the ST617 phylogeny with  $\beta$ -lactamase gene carriage shows a high prevalence of the *bla*<sub>CTX-M-15</sub> ESBL gene in characterised isolates (Fig 2A). Annotation of the ST167 phylogeny with  $\beta$ -lactamase gene carriage (Fig 2B) shows a pattern of resistance gene carriage, with multiple independent acquisitions of carbapenemase across the phylogeny including *bla*<sub>NDM-1</sub>, *bla*<sub>NDM-5</sub>, *bla*<sub>NDM-7</sub>, *bla*<sub>OXA-181</sub>, and *bla*<sub>KPC-3</sub>. For both phylogenies there is clear evidence of isolation of strains from across the globe.

#### **Evolutionary genomic analysis correlates switches in LPS gene content with the emergence of the ST167/ST617 lineage**

Both ST167 and ST617 are single locus variants of the ST10 lineage of *E. coli*. ST10 is the most abundant lineage of *E. coli* represented in the Enterobase database and contains isolates ranging from drug susceptible environmental and human commensal strains, to multi-drug resistant strains isolated from human clinical UTI and bacteraemia infections. We selected 256 ST10 genomes from Enterobase (Table S4) to represent the known spectrum of ST10 diversity present in the database, and merged this data set with our publically available ST167/ST617 genome data set to create a larger ST10 complex phylogeny (Fig S1). The resulting phylogeny shows that ST167 and ST617 are sister clades with respect to ST10, with ST617 emerging as a nested clade from a single outlying ST167 genome, though the distance between ST167 and ST617 is around 18,000 SNPs.

Given the phylogenetic pattern of ST167 and ST617 with respect to ST10, we sought to determine if their emergence from ST10 is associated with defined evolutionary events. We used a combined GWAS approach to compare the ST167/617 genomes

235 with ST10, using both SEER and SCOARY analysis of a pangenome matrix. Only  
236 loci considered to be significantly associated with one lineage over the other by both  
237 methods were further investigated (Dataset S1). Most striking was the absence of  
238 the *wzzB* gene and *wca* biosynthetic cluster in ST167/ST617 whilst the majority of  
239 the ST10 genomes contained both (Figure S2). These genes are involved in LPS  
240 biosynthesis with *wzzB* being the master controller of O antigen chain length in the  
241 *wzx/wzy* pathway, whilst *wca* genes are responsible for colonic acid biosynthesis.<sup>31</sup>  
242 In silico *E. coli* serotyping<sup>32</sup> established that ST167 and ST617 demonstrate the  
243 exact same O antigenic type (O32novel) with similarity also seen in H antigen type  
244 (H9 or H10) (Figure S2), whilst the SerotypeFinder database identified the strains as  
245 O89.

246 Our combined GWAS analysis also identified another ~90 CDS which were present  
247 across the entire data set, but which had distinct alleles in the ST167/ST617  
248 genomes compared to those in ST10 (Fig 3, Dataset S2). Many of these CDS  
249 encode dehydrogenase enzymes involved in anaerobic metabolism, or are part of  
250 the *cob/pdu/eut* operons known to be involved in anaerobic respiration during  
251 intestinal inflammation.<sup>33</sup> This would appear to suggest differential evolutionary  
252 events in key genes involved in anaerobic metabolism in the formation of the  
253 ST167/ST617 lineage. Also present were unique alleles in core CDS involved in acid  
254 and bile salt tolerance, and a number of fimbrial-like proteins. In conjunction these  
255 data would suggest differential evolutionary forces acting on loci involved in  
256 mammalian colonisation in ST167/617 in comparison to ST10. Furthermore a  
257 combined SEER and Piggy approach identified unique sequences in 17 intergenic  
258 regions (IGRs) upstream of core CDS in ST167/617 that were distinct from ST10,

including IGRs upstream of anaerobic metabolic loci also present in the SEER/SCOARY analysis (Dataset S1).

## Discussion

Our data presented here provide a comprehensive genomic analysis of two lineages of carbapenem resistant *E. coli* infecting multiple patients within the ICU of West China hospital. Both these lineages, ST167 and ST617, are members of the larger ST10 complex of *E. coli*, which is ubiquitously found in environmental, human clinical, and mammalian intestinal commensal sampling. Our analysis is the first genome level characterisation of strains belonging to ST167 or ST617, despite a number of single site reports of clinical infections with both lineages existing in the literature.

Our analysis shows that the diversity which accumulates in the genome of the ST167 isolates during the course of the investigation is not mirrored by diversity in the plasmid carrying the *bla<sub>NDM-5</sub>* gene. Only 1 SNP difference existed between the sequence of this plasmid in the ST167 isolates, and only 2 SNPs difference between the ST167 and ST617 isolates. As a result it is impossible to tell if the IncX3 plasmid associated with dissemination of *bla<sub>NDM-5</sub>* in China<sup>34</sup> was transferred between ST167 and ST617 in the hospital, or if the plasmid is highly stable with only deleterious mutations occurring and quickly purged from the population. Clearly there is a need for more thorough and detailed analysis of various resistance plasmids within and between hospitals, such as was done recently for NDM-1 plasmids in Latin America.

<sup>35</sup>

The lack of appropriately designed isolate collection and sequencing strategy means it is impossible to conduct any form of genomic epidemiological analyses of these *E. coli* lineages beyond our Chinese investigation. However the ready availability of a

284 large number of good-quality, curated genome assemblies in the Enterobase  
285 genome database do allow us to delve deeper into the evolutionary history of *E. coli*  
286 ST167 and ST617. Whilst data generated and uploaded to Enterobase is prone to a  
287 bias towards clinical MDR strains, it is still clear that ESBL and carbapenem resistant  
288 strains of both these lineages have been isolated from across the world over the past  
289 20 or so years (Tables S1 and S2).

290 Comparative genomic analysis and GWAS for traits specific to ST167 and ST617  
291 compared to ST10 also support emergence along a shared evolutionary branch. Key  
292 among these is the complete loss of the *wca* operon encoding colanic acid  
293 biosynthesis in the LPS biosynthesis pathway. The majority of *E. coli* produce their  
294 LPS utilising the O-unit translocation pathway encoded for by *wzx* and *wzy*.<sup>31</sup> This  
295 method utilises glycosyltransferases to assemble the O antigen in units at the  
296 cytoplasmic membrane. These units are then translocated by Wzx and polymerized  
297 by Wzy until the O antigen chain length is reached. This mechanism is utilised by the  
298 majority of the ST10 isolates, however genomic analysis shows that ST167 and  
299 ST617 utilise an alternative *wzm/wzt* ATP transporter pathway. This biosynthetic  
300 pathway assembles the entire O-antigen on the cytoplasmic face before Wzt  
301 transports the O-chain across,<sup>31</sup> resulting in an O-antigen with truncated chain  
302 length. O-antigen chain length plays a major role in pathogenicity of Gram negative  
303 organisms, and it has been demonstrated that loss of long O-antigen chains in  
304 *Salmonella* optimizes immune evasion and allows successful colonisation.<sup>36</sup>

305 Alongside the LPS genetic changes, we also observed unique alleles of anaerobic  
306 metabolism genes and genes potentially involved in host colonisation in ST167/617  
307 compared to ST10. Recent modelling data has shown that any factor influencing the

ability of a bacterium to colonise a host will also influence its likelihood of evolving antimicrobial resistance.<sup>37</sup>

## **Conclusions**

We provide data for the first ever, single hospital genomic analysis of clinical isolates of carbapenem resistant *E. coli* belonging to the ST167/617 lineage. Our data presented here provide evidence for evolutionary events that would affect microbial interaction with a mammalian host underpinning the emergence of the ST167/617 lineage from ST10. There is also evidence for lineage specific alterations in intergenic regions in ST167/617, a phenomenon which has already been described as underpinning the emergence of MDR plasmid-containing *E. coli* ST131 strains.<sup>28</sup> Clearly there is now a need for a fully designed genomic epidemiological investigation of lineages of *E. coli* associated with carriage of carbapenem resistance plasmids arising from the ST10 clade, both in China and internationally. Such a study will fully inform us of any potential parallelism in the evolution of MDR lineages of *E. coli*, and of the true nature and scope of their prevalence and global dissemination.

## **Funding section**

This work was funded by a Royal Society Newton Advanced Fellowship project (NA150363) and a grant from the National Natural Science Foundation of China (project no. 8151101182) awarded to ZZ and AM. SF was funded by the Wellcome Antimicrobial Resistance doctoral training project at UoB, and CC by the Wellcome MIDAS doctoral training program at UoB.

## **Transparency declaration**

None to declare

## References

1. de Kraker MEA, Jarlier V, Monen JCM, et al. The changing epidemiology of bacteraemias in Europe: trends from the European Antimicrobial Resistance Surveillance System. *Clin. Microbiol. Infect.* 2013;**19**:860–8.
2. Mathers AJ, Peirano G, Pitout JDD. The role of epidemic resistance plasmids and international high-risk clones in the spread of multidrug-resistant Enterobacteriaceae. *Clin. Microbiol. Rev.* 2015;**28**:565–91.
3. Alhashash F, Weston V, Diggle M, McNally A. Multidrug-Resistant *Escherichia coli* Bacteremia. *Emerg. Infect. Dis.* 2013;**19**:1699–701.
4. Croxall G, Hale J, Weston V, et al. Molecular epidemiology of extraintestinal pathogenic *Escherichia coli* isolates from a regional cohort of elderly patients highlights the prevalence of ST131 strains with increased antimicrobial resistance in both community and hospital care settings. *J. Antimicrob. Chemother.* 2011;**66**:2501–8.
5. Livermore DM, Hawkey PM. CTX-M: changing the face of ESBLs in the UK. *J. Antimicrob. Chemother.* 2005;**56**:451–4.
6. Feng Y, Yang P, Xie Y, et al. *Escherichia coli* of sequence type 3835 carrying blaNDM-1, blaCTX-M-15, blaCMY-42 and blaSHV-12. *Sci. Rep.* 2015;**5**:12275.
7. Zhang L, Xue W, Meng D. First report of New Delhi metallo-beta-lactamase 5 (NDM-5)-producing *Escherichia coli* from blood cultures of three leukemia patients. *Int. J. Infect. Dis.* 2016;**42**:45–6.
8. Cuzon G, Bonnin RA, Nordmann P. First identification of novel NDM carbapenemase, NDM-7, in *Escherichia coli* in France. *PLoS One.* 2013;**8**:e61322.
9. Zheng B, Dong H, Xu H, et al. Coexistence of MCR-1 and NDM-1 in Clinical *Escherichia coli* Isolates. *Clin. Infect. Dis.* 2016;**63**:1393–5.



358 10. Leflon-Guibout V, Blanco J, Amaqdouf K, et al. Absence of CTX-M Enzymes but  
359 High Prevalence of Clones, Including Clone ST131, among Fecal *Escherichia coli*  
360 Isolates from Healthy Subjects Living in the Area of Paris, France. *J. Clin. Microbiol.*  
361 2008;**46**:3900–5.

362 11. Kohler C-D, Dobrindt U. What defines extraintestinal pathogenic *Escherichia*  
363 *coli*? *Int. J. Med. Microbiol.* 2011;**301**:642–7.

364 12. Cerqueira GC, Earl AM, Ernst CM, et al. Multi-institute analysis of carbapenem  
365 resistance reveals remarkable diversity, unexplained mechanisms, and limited clonal  
366 outbreaks. *Proc. Natl. Acad. Sci.* 2017;**114**:1135–40.

367 13. Bankevich A, Nurk S, Antipov D, et al. SPAdes: a new genome assembly  
368 algorithm and its applications to single-cell sequencing. *J. Comput. Biol.*  
369 2012;**19**:455–77.

370 14. Seemann T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics.*  
371 2014;**30**:2068–9.

372 15. Larsen M V, Cosentino S, Rasmussen S, et al. Multilocus sequence typing of  
373 total-genome-sequenced bacteria. *J. Clin. Microbiol.* 2012;**50**:1355–61.

374 16. Treangen TJ, Ondov BD, Koren S, Phillippy AM. The Harvest suite for rapid  
375 core-genome alignment and visualization of thousands of intraspecific microbial  
376 genomes. *Genome Biol.* 2014;**15**:524.

377 17. Loman NJ, Quinlan AR. Poretools: a toolkit for analyzing nanopore sequence  
378 data. *Bioinformatics.* 2014;**30**:3399–401.

379 18. Koren S, Walenz BP, Berlin K, et al. Canu: scalable and accurate long-read  
380 assembly via adaptive k-mer weighting and repeat separation. *Genome Res.*  
381 2017;**27**:722–36.

382 19. Walker BJ, Abeel T, Shea T, et al. Pilon: an integrated tool for comprehensive

383 microbial variant detection and genome assembly improvement. *PLoS One*.  
384 2014;**9**:e112963.

385 20. Page AJ, Cummins CA, Hunt M, et al. Roary: rapid large-scale prokaryote pan  
386 genome analysis. *Bioinformatics*. 2015;**31**:3691–3.

387 21. Stamatakis A Ludwig T MH. RAXML-III: a fast program for maximum likelihood-  
388 based inference of large phylogenetic trees. *Bioinformatics*. 2005;**21**:456.

389 22. Letunic I, Bork P. Interactive Tree Of Life v2: online annotation and display of  
390 phylogenetic trees made easy. *Nucleic Acids Res*. 2011;**39**:W475-8.

391 23. Brynildsrud O, Bohlin J, Scheffer L, Eldholm V. Rapid scoring of genes in  
392 microbial pan-genome-wide association studies with Scoary. *Genome Biol*.  
393 2016;**17**:238.

394 24. Lees JA, Vehkala M, Valimaki N, et al. Sequence element enrichment analysis to  
395 determine the genetic basis of bacterial phenotypes. *Nat. Commun*. 2016;**7**:12797.

396 25. Inouye M, Conway TC, Zobel J, Holt KE. Short read sequence typing (SRST):  
397 multi-locus sequence types from short reads. *BMC Genomics*. 2012;**13**:338.

398 26. Joensen KG, Tetzschner AMM, Iguchi A, et al. Rapid and Easy In Silico  
399 Serotyping of *Escherichia coli* Isolates by Use of Whole-Genome Sequencing Data.  
400 *J. Clin. Microbiol*. 2015;**53**:2410–26.

401 27. Thorpe HA, Bayliss SC, Sheppard SK, Feil EJ. Piggy: A Rapid, Large-Scale Pan-  
402 Genome Analysis Tool for Intergenic Regions in Bacteria. *Gigascience*. 2018;**7**:1-11.

403 28. McNally A, Oren Y, Kelly D, et al. Combined Analysis of Variation in Core,  
404 Accessory and Regulatory Genome Regions Provides a Super-Resolution View into  
405 the Evolution of Bacterial Populations. *PLoS Genet*. 2016;**12**:e1006280.

406 29. Köser CU, Holden MT, Ellington MJ, et al. Rapid whole-genome sequencing for  
407 investigation of a neonatal MRSA outbreak. *N. Engl. J. Med*. 2012;**366**:2267–75.

408 30. Quick J, Cumley N, Wearn CM, et al. Seeking the source of *Pseudomonas*  
409 *aeruginosa* infections in a recently opened hospital: an observational study using  
410 whole-genome sequencing. *BMJ Open*. 2014;**4**:e006278.

411 31. Iguchi A, Iyoda S, Kikuchi T, et al. A complete view of the genetic diversity of the  
412 *Escherichia coli* O-antigen biosynthesis gene cluster. *DNA Res*. 2015;**22**:101–7.

413 32. Ingle DJ, Valcanis M, Kuzevski A, et al. In silico serotyping of *E. coli* from short  
414 read data identifies limited novel O-loci but extensive diversity of O:H serotype  
415 combinations within and between pathogenic lineages. *Microb. genomics*.  
416 2016;**2**:e000064.

417 33. McNally A, Thomson NR, Reuter S, Wren BW. “Add, stir and reduce”: *Yersinia*  
418 *spp.* as model bacteria for pathogen evolution. *Nat. Rev. Microbiol*. 2016;**14**:177–90.

419 34. Yang P, Xie Y, Feng P, Zong Z. blaNDM-5 carried by an IncX3 plasmid in  
420 *Escherichia coli* sequence type 167. *Antimicrob. Agents Chemother*. 2014;**58**:7548–  
421 52.

422 35. Marquez-Ortiz RA, Haggerty L, Olarte N, et al. Genomic Epidemiology of NDM-1-  
423 Encoding Plasmids in Latin American Clinical Isolates Reveals Insights into the  
424 Evolution of Multidrug Resistance. *Genome Biol. Evol*. 2017;**9**:1725–41.

425 36. Crawford RW, Wangdi T, Spees AM, et al. Loss of very-long O-antigen chains  
426 optimizes capsule-mediated immune evasion by *Salmonella enterica* serovar Typhi.  
427 *MBio*. 2013;**4**: e00232-13.

428 37. Lehtinen S, Blanquart F, Croucher NJ, et al. Evolution of antibiotic resistance is  
429 linked to any genetic mechanism affecting bacterial duration of carriage. *Proc. Natl.*  
430 *Acad. Sci*. 2017;**114**:1075–80.

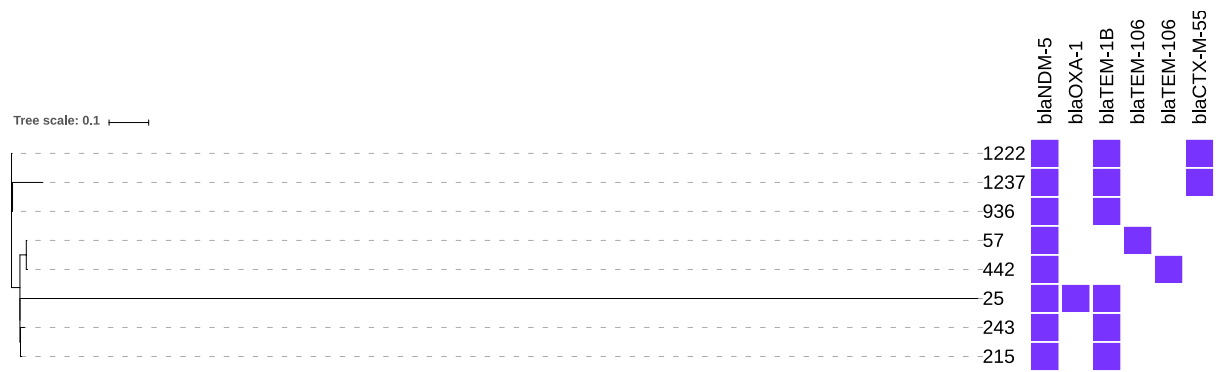
431

432

Figure 1: Maximum likelihood phylogenetic tree of *E. coli* ST167 strains isolated from the ICU of West China hospital. The phylogeny is inferred from a SNP alignment obtained by mapping raw data against a Minlon/Illumina hybrid complete assembly of isolate 1237. The annotation denotes the presence of ESBL and CPE associated  $\beta$ -lactamases as determined by Abricate.

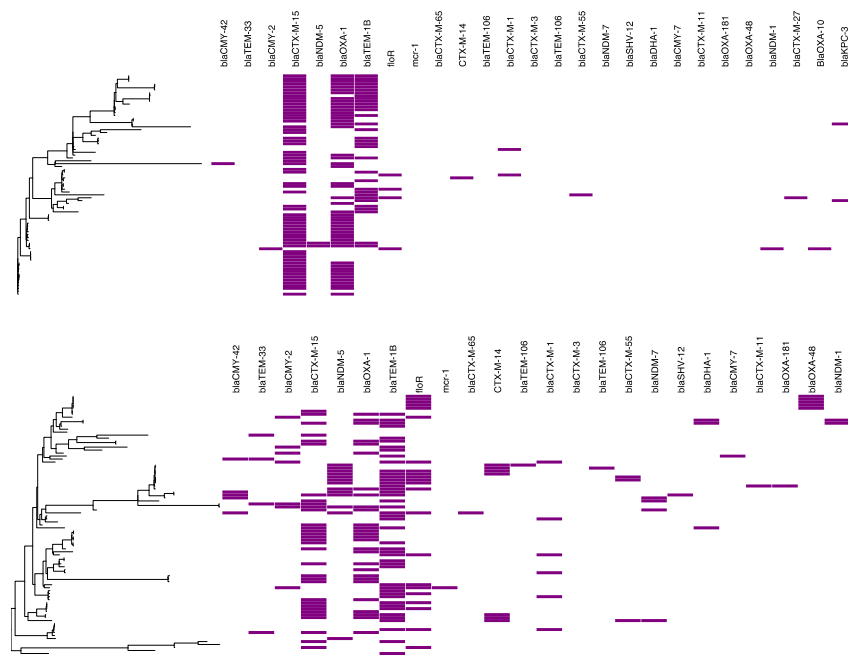
Figure 2: Maximum likelihood phylogenetic trees of a global collection of (A) ST617 and (B) ST167 strains. The phylogeny is inferred from an alignment of concatenated core CDS sequences as determined by Roary, and is mid-point rooted. The annotation denotes the presence of ESBL and CPE associated  $\beta$ -lactamases as determined by Abricate.

Figure 3: Manhattan skyline plot showing position of kmers identified by GWAS analysis as being significantly associated with ST167/617 compared to ST10. The x axis indicates the position on the WCHec1237 complete genome assembly, whilst the Y axis indicates the numbers of statistically significant kmers mapping at that position. Hits indicated in red are either intergenic regions (labelled IGR) identified as being unique by both Piggy and SEER analysis, or anaerobic metabolism loci identified as significantly different by both SEER and Scoary.

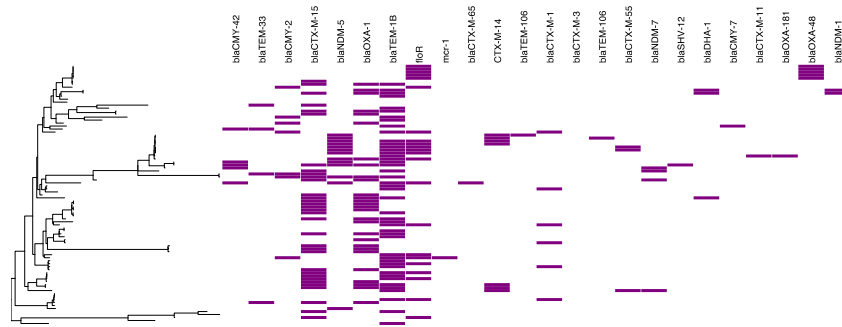


SNP Distance Matrix	1222	215	243	25	442	57	936	1237
1222	0	30	39	725	45	43	5	74
215	30	0	15	705	25	23	31	100
243	39	15	0	714	34	32	40	113
25	725	705	714	0	720	718	726	799
442	45	25	34	720	0	6	46	113
57	43	23	32	718	6	0	44	117
936	5	31	40	726	46	44	0	73
1237	74	100	113	799	113	117	73	0

Figure 1

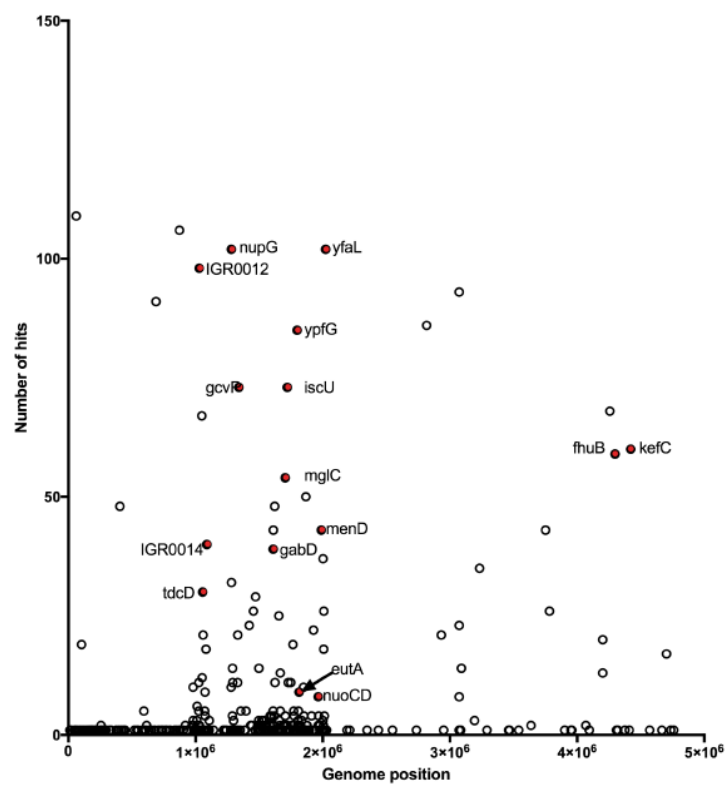


**A**



**B**

Figure 2



460

461 Figure 3

462